

IMPROVING NEWTON'S METHOD PERFORMANCE BY PARAMETRIZATION: THE CASE OF RICHARDS EQUATION

KONSTANTIN BRENNER AND CLÉMENT CANCÈS

ABSTRACT. The nonlinear systems obtained by discretizing degenerate parabolic equations may be hard to solve, especially with Newton's method. In this paper, we apply to Richards equation a strategy that consists in defining a new primary unknown for the continuous equation in order to stabilize Newton's method by parametrizing the graph linking the pressure and the saturation. The resulting form of Richards equation is then discretized thanks to a monotone Finite Volume scheme. We prove the well-posedness of the numerical scheme. Then we show under appropriate non-degeneracy conditions on the parametrization that Newton's method converges locally and quadratically. Finally, we provide numerical evidences of the efficiency of our approach.

Keywords. Richards equation, Finite Volumes, Newton's method, parametrization

AMS subjects classification. 65M22, 65M08, 76S05

1. INTRODUCTION

1.1. Motivations and presentation of the Richards equation. Solving numerically some nonlinear partial differential equations, for example by using finite elements or finite volumes, often amounts to the resolution of some nonlinear system of equations of the form:

$$(1) \quad \text{Find } \mathbf{u} \in \mathbb{R}^N \text{ such that } \mathcal{R}(\mathbf{u}) = \mathbf{0}_{\mathbb{R}^N},$$

where $N \in \mathbb{N}^*$ is the number of degrees of freedom and can be large. One of the most popular method for solving the systems of the form (1) is the celebrated Newton-Raphson method. If this iterative procedure converges, then its limit is necessarily a solution to (1). However, making the Newton method converge is sometimes difficult and might require a great expertise. Nonlinear preconditioning technics have been recently developed in improve the performance of the Newton's method, see for instance [18, 7].

Complex multiphase or unsaturated porous media flows are often modeled thanks to degenerate parabolic problems. We refer to [3] for an extensive discussion about models of porous media flows. For such degenerate problems, making Newton's method converge is often very difficult. This led to the development of several strategies to optimize the convergence properties, like for instance the so-called continuation-Newton method [45], or trust region based solvers [44]. An alternative approach consist in solving (1) thanks to a robust fixed point procedure with linear convergence speed rather than with the quadratic Newton's method (see for instance [46, 30, 31, 40]). Comparisons between the fixed point and the Newton's

This work was supported by the GeoPor project funded by the French National Research Agency (ANR) with the grant ANR-13-JS01-0007-01 (project GEOPOR) .

strategies are presented for instance in [33, 4] (see also [42]). Combinations of both technics (perform few fixed points iterations before running Newton's algorithm) was for instance performed in [35].

Our strategy consists in reformulating the problem before applying Newton's method. The reformulation consists in changing the primary variable in order to improve the behavior of Newton's method. We apply this strategy to the so called *Richards equation* [43, 3] modeling the unsaturated flow of water within a porous medium. Extension to more complex models of porous media flows will be the purpose of the forthcoming contribution [5].

Denote by Ω some open subset of \mathbb{R}^d ($d \leq 3$) representing the porous medium (in the sequel, Ω will be supposed to be polyhedral for meshing purpose), by $T > 0$ a finite time horizon, and by $Q := \Omega \times (0, T)$ the corresponding space-time cylinder. We are interested in finding a saturation profile $\bar{s} : Q \rightarrow [0, 1]$ and a water pressure $\bar{p} : Q \rightarrow \mathbb{R}$ such that

$$(2) \quad \partial_t \bar{s} - \nabla \cdot (\lambda(\bar{s}) (\nabla \bar{p} - \mathbf{g})) = 0,$$

where the mobility function $\lambda : [0, 1] \rightarrow \mathbb{R}_+$ is a nondecreasing \mathcal{C}^2 function that satisfies $\lambda(s \leq 0) = 0$ and $\lambda(s > 0) > 0$, and where $\mathbf{g} \in \mathbb{R}^d$ stands for the gravity vector. In order to ease the reading, we have set the porosity equal to 1 in (2) and neglected the residual saturation. The pressure and the water content are supposed to be linked by some monotone relation

$$(3) \quad \bar{s} = S(\bar{p}), \quad \text{a.e. in } Q$$

where S is a non-decreasing function from \mathbb{R} to $[0, 1]$. In what follows, we assume that $S(p) = 1$ for all $p \geq 0$, that corresponds to assuming that the porous medium is water wet, and that $S \in L^1(\mathbb{R}_-)$, implying in particular that $\lim_{p \rightarrow -\infty} S(p) = 0$. As a consequence of the Lipschitz continuity of λ and of the integrability of S on \mathbb{R}_- , one has

$$(4) \quad \lambda(S) \in L^1(\mathbb{R}_-).$$

We denote by $p_\star = \sup\{p \mid S(p) = 0\}$, with the convention that $p_\star = -\infty$ if $\{p \in \mathbb{R} \mid S(p) = 0\} = \emptyset$.

Typical behaviors of λ and S are depicted in Figure 1.

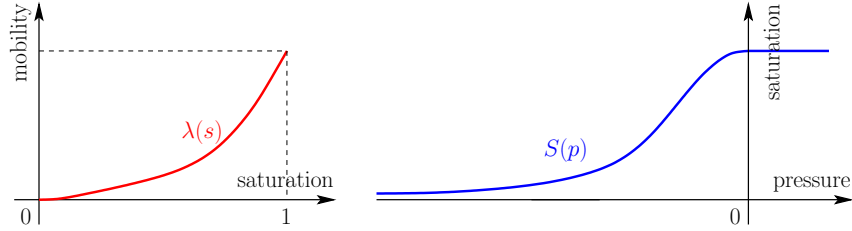


FIGURE 1. The mobility function $\lambda : [0, 1] \rightarrow \mathbb{R}_+$ is increasing and satisfies $\lambda(0) = 0$. The saturation function $S : \mathbb{R} \rightarrow [0, 1]$ is non-decreasing, constant equal to 1 on \mathbb{R}_+ and increasing on $(p_\star, 0)$ for some $p_\star \in [-\infty, 0)$.

Remark 1.1. *In the case where the domain Ω and the finite time T are large, a classical hyperbolic scaling consisting in replacing (\mathbf{x}, t) by $(\mathbf{x}/\epsilon, t/\epsilon)$ leads to the problem*

$$\partial_t S^\epsilon(\bar{p}) - \nabla \cdot (\lambda(S^\epsilon(\bar{p})) (\nabla \bar{p} - \mathbf{g})) = 0,$$

where the function S^ϵ is deduced from S by

$$S^\epsilon(p) = S(p/\epsilon), \quad \forall p \in \mathbb{R}.$$

Letting ϵ tend to 0 leads the maximal monotone capillary pressure graph

$$S^0(p) = \text{sign}_+(p) = \begin{cases} 0 & \text{if } p < 0, \\ [0, 1] & \text{if } p = 0, \\ 1 & \text{if } p > 0. \end{cases}$$

The Richards equation then degenerates into a hyperbolic-elliptic problem. Our purpose can be extended to the degenerate case even though the graph S^0 lack regularity thanks to the so-called semi-smooth Newton method (cf. [41]). Moreover, because of the hyperbolic degeneracy, additional entropy criterions à la Carrillo [13] are required in order to characterize the relevant solution. In order to simplify our purpose as much as possible, we focus on the case $\epsilon > 0$.

The problem (2)–(3) is complemented by the initial condition

$$(5) \quad \bar{s}|_{t=0} = s_0 \in L^\infty(\Omega; [0, 1]),$$

and by Dirichlet boundary conditions on the pressure:

$$(6) \quad \bar{p}|_{\partial\Omega \times (0, T)} = p_D.$$

The regularity requirements on the boundary condition will be specified later on.

At least from a mathematical point of view, it is natural to solve (2) by choosing p or the Kirchhoff transform u (to be defined later on at (8)) as a primary unknown then to deduce $s = S(p) = \tilde{S}(u)$. However, this approach lacks efficiency when one aims to solve the Richards equation numerically, especially for dry media, i.e., when the saturation s is close to 0. In this latter situation, it turns out that the Newton methods encounters difficulties to converge for solving the nonlinear system obtained thanks to standard implicit numerical methods (say $\mathbb{P}1$ -Finite Elements [33], mixed finite elements [4], or Finite Volumes [25, 27]). A better choice as a primary unknown in the dry regions is the saturation s , p or u being computed thanks to the inverse function of S or \tilde{S} respectively. But choosing the saturation s as the primary variable yields difficulties in the saturated regions, i.e., where $s = 1$. Hence, a classical approach for solving numerically the Richards equation consists in applying the so-called *variable switch*, that consists in changing the primary variable following the physical configuration (see, e.g., [17]).

1.2. Monotone parametrization of the graph. The main feature of our contribution consists in parametrizing the graph S in order to stabilize the Newton algorithm without implementing the possibly complex variable switch procedure. This procedure is inspired from the one proposed by J. Carrillo [14] (see also [36] for numerical issues) to deal with hyperbolic scalar conservation laws with discontinuous flux w.r.t. the unknown. Let us introduce two continuously differentiable nondecreasing functions

$$s : (\tau_\star, \infty) \rightarrow [0, 1] \quad \text{and} \quad p : (\tau_\star, \infty) \rightarrow (-\infty, \infty),$$

where $\tau_\star < 0$ may be equal to $-\infty$, such that $p(0) = 0$ and

$$\bar{s} \in S(\bar{p}) \quad \Leftrightarrow \quad \text{there exists } \tau \geq \tau_\star \text{ s.t. } \bar{s} = s(\tau) \text{ and } \bar{p} = p(\tau).$$

This enforces in particular that $\lim_{\tau \rightarrow \tau_\star} s(\tau) = 0$ and $\lim_{\tau \rightarrow \tau_\star} p(\tau) = p_\star$. In the case where $\tau_\star > -\infty$, the functions s and p are then continuously extended into constants on $(-\infty, \tau_\star)$. It is assumed that the parametrization function s satisfies

$$(7) \quad 1 - s \in L^1(\mathbb{R}_+) \text{ and } s \in L^1(\mathbb{R}_-).$$

The Kirchhoff transform $u : [\tau_\star, +\infty) \rightarrow \mathbb{R}$ is defined by

$$(8) \quad u(\tau) = \int_0^\tau \lambda(s(a))p'(a)da, \quad \forall \tau \geq \tau_\star.$$

It follows from the integrability property (4) that

$$(9) \quad u_\star := \lim_{\tau \searrow \tau_\star} u(\tau) \text{ is finite.}$$

For technical reasons, the function u is artificially extended into a continuous onto function from \mathbb{R} to \mathbb{R} by setting

$$(10) \quad u(\tau) = \tau - \tau_\star + u_\star, \quad \forall \tau < \tau_\star.$$

However, as it will appear later on (cf. Lemma 2.7), the choice of this extension has no influence on the result.

It is assumed throughout this paper that the parametrization is not degenerated, i.e.,

$$s'(\tau) + p'(\tau) > 0 \quad \text{for a.e. } \tau \geq \tau_\star,$$

or equivalently

$$s'(\tau) + u'(\tau) > 0 \quad \text{for a.e. } \tau \in \mathbb{R}$$

since $\lambda(s(\tau)) > 0$ for all $\tau > \tau_\star$. Since S is an absolutely continuous function, this implies in particular that $p' > 0$ a.e. in \mathbb{R}_+ .

Such a parametrization of the graph S always exists but is not unique. For instance, one can choose the parametrizations defined by $p(\tau) = \tau$ or $u(\tau) = \tau$. As it will appear in the analysis carried out in the core of the paper, a convenient parametrization should satisfy: there exist $\alpha_\star > 0$ and $\alpha^\star \geq \alpha_\star$ such that

$$(11) \quad \alpha_\star \leq \max(s'(\tau), u'(\tau)) \leq \alpha^\star, \quad \forall \tau \in \mathbb{R}.$$

Thus, we assume that (11) holds for the analysis. This ensures in particular that the functions s and u are Lipschitz continuous:

$$(12) \quad \|s'\|_\infty \leq \alpha^\star, \quad \|u'\|_\infty \leq \alpha^\star.$$

For technical reasons that will appear in the analysis, we also assume that there exists $C > 0$ such that

$$(13) \quad \tau \leq C(u(\tau) + 1), \quad \forall \tau \geq 0.$$

It is also assumed that

$$(14) \quad \liminf_{\tau \searrow \tau_\star} p'(\tau) > 0.$$

This assumption is very naturally satisfied for any non-degenerate parametrization in the sense of (11) of a reasonable function S , but unphysical counterexamples can be designed, enforcing us to set (14) as an assumption.

The function s from $[\tau_*, 0]$ to $[0, 1]$ is nondecreasing and onto. Therefore, one can define the function $s^{-1} : [0, 1] \rightarrow [\tau_*, 0]$ by

$$(15) \quad s^{-1}(a) = \min\{x \geq 0 \mid s(x) = a\}, \quad \forall a \in [0, 1].$$

This allows to define an initial data τ_0 as $\tau_0 = s^{-1}(s_0)$ such that $s(\tau_0) = s_0$.

Choosing τ as the primary variable leads to the following doubly degenerate parabolic equation

$$\partial_t s(\tau) - \nabla \cdot \left(\lambda(s(\tau)) (\nabla p(\tau) - \mathbf{g}) \right) = 0 \quad \text{in } Q.$$

This equation turns to

$$(16) \quad \partial_t s(\tau) + \nabla \cdot \left(\lambda(s(\tau)) \mathbf{g} - \nabla u(\tau) \right) = 0 \quad \text{in } Q,$$

at least if $\tau \geq \tau_*$ (this will be ensured, cf. Theorem 1.3). It is relevant to impose the boundary condition

$$(17) \quad \tau|_{\partial\Omega \times (0, T)} = p^{-1}(p_D) =: \tau_D \geq \tau_*.$$

as a counterpart of (6). It is finally assumed that τ_D can be extended to the whole $\Omega \times (0, T)$ in a way such that

$$(18) \quad \tau_D \in C^1(\overline{Q}), \quad \text{with } \tau \geq \tau_*.$$

The regularity required on τ_D is not optimal and can be relaxed. However, the treatment of the boundary condition is not central in our purpose, hence we stick to (18)

Definition 1.2. *A measurable function $\tau : Q \rightarrow \mathbb{R}$ is said to be a weak solution to the problem (16), (5), (17) if $u(\tau) - u(\tau_D) \in L^2((0, T); H_0^1(\Omega))$, if $\partial_t s(\tau) \in L^2((0, T); H^{-1}(\Omega))$, and if, for all $\varphi \in C_c^\infty(\Omega \times [0, T]; \mathbb{R})$, one has*

$$\iint_Q s(\tau) \partial_t \varphi \, d\mathbf{x} dt + \int_\Omega s_0 \varphi(\mathbf{x}, 0) \, d\mathbf{x} + \iint_Q \left(\lambda(s(\tau)) \mathbf{g} - \nabla u(\tau) \right) \cdot \nabla \varphi \, d\mathbf{x} dt = 0.$$

The following statement summarizes known results about the weak solutions.

Theorem 1.3. *There exists a unique weak solution $\tau : \Omega \rightarrow \mathbb{R}$ to the problem (16), (5), (17) in the sense of Definition 1.2. Moreover, $\tau \geq \tau_*$ a.e. in Q , $s(\tau) \in C([0, T]; L^p(\Omega))$ for all $p \in [1, \infty)$, and, given two solutions $\tau, \hat{\tau}$ corresponding to two initial data s_0 and \hat{s}_0 , we have*

$$(19) \quad \int_\Omega (s(\tau(\mathbf{x}, t)) - s(\hat{\tau}(\mathbf{x}, t)))^\pm \, d\mathbf{x} \leq \int_\Omega (s_0(\mathbf{x}) - \hat{s}_0(\mathbf{x}))^\pm \, d\mathbf{x}, \quad \forall t \in [0, T].$$

Existence of weak solutions have been proved by Alt and Luckhaus in their seminal paper [1]. We refer to [39] (see also [12, 28]) for extended details on the uniqueness proof and on the comparison principle (19). The time continuity of the saturation can be proved as in [8].

1.3. Outline of the paper. In §2, we present an implicit monotone Finite Volume scheme [23] designed for approximating the entropy solution τ of (5) and (16)–(17). First, we describe in §2.1 how the domains Ω (and then Q) has to be meshed. In particular, the mesh has to fulfill the so-called *orthogonality condition* so that the diffusion fluxes can be discretized using a simple two-point flux approximation [22].

The Finite Volume scheme is described in §2.2. This scheme yields a nonlinear system of equations

$$(20) \quad \mathcal{F}_n(\boldsymbol{\tau}^n) = \mathbf{0}, \quad \forall n \in \{1, \dots, N\}$$

to be solved at each time step. The existence and the uniqueness of the solution $\boldsymbol{\tau}^n$ to the nonlinear system (20) is proved at §2.3.

Once we know that the scheme (20) admits one unique solution $\boldsymbol{\tau}^n$, we discuss its effective computation thanks to Newton's method in §3. It is in particular proved in §3.2 that the jacobian matrix is uniformly non-degenerate, so that we can use Newton-Kantorovich theorem to claim the convergence of Newton's method. The quantification of the error linked to the inexact resolution of the nonlinear system is performed in §3.3. Finally, some numerical results are presented in §4 in order to illustrate the efficiency of our approach.

2. THE FINITE VOLUME SCHEME

2.1. Discretization of Q . In this work, we only consider cylindrical discretizations of Q that consist in discretizing space and time separately.

2.1.1. Admissible mesh of Ω . The approximation of the diffusive fluxes we propose relies on the so-called *two-point flux approximation*. This approximation is consistent if the problem is isotropic and if the mesh satisfies the so-called orthogonality condition (see e.g. [22]).

Definition 2.1 (admissible mesh of Ω). *An admissible mesh $(\mathcal{T}, \mathcal{E}, (\mathbf{x}_K)_{K \in \mathcal{T}})$ of Ω is given by a set \mathcal{T} of disjointed open bounded convex subsets of Ω called control volumes, a family \mathcal{E} of subsets of $\bar{\Omega}$ called edges contained in hyperplanes of \mathbb{R}^d with strictly positive measure, and a family of points $(\mathbf{x}_K)_{K \in \mathcal{T}}$ (the so-called cell centers). It is assumed that the mesh integrates the whole Ω , i.e., $\bigcup_{K \in \mathcal{T}} \bar{K} = \bar{\Omega}$. The boundary of the control volumes are made of edges, i.e., for all $K \in \mathcal{T}$, there exists a subset \mathcal{E}_K of \mathcal{E} such that $\partial K = \bigcup_{\sigma \in \mathcal{E}_K} \bar{\sigma}$. Furthermore, $\mathcal{E} = \bigcup_{K \in \mathcal{T}} \mathcal{E}_K$. For any $(K, L) \in \mathcal{T}^2$ with $K \neq L$, either the $(d-1)$ -dimensional Lebesgue measure of $\bar{K} \cap \bar{L}$ is 0, or $\bar{K} \cap \bar{L} = \bar{\sigma}$ for some $\sigma \in \mathcal{E}$. In the latter case, we write $\sigma = K|L$. We denote by $\mathcal{E}_{\text{int}} = \{\sigma \in \mathcal{E}, \exists (K, L) \in \mathcal{T}^2 \sigma = K|L\}$ the set of the internal edges, and by $\mathcal{E}_{\text{ext}} = \{\sigma \in \mathcal{E}, \sigma \subset \partial\Omega\}$, $\mathcal{E}_{K, \text{ext}} = \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$ of the boundary edges. Finally, the family of points $(\mathbf{x}_K)_{K \in \mathcal{T}}$ is such that $\mathbf{x}_K \in K$ (for all $K \in \mathcal{T}$) and, if $\sigma = K|L$, it is assumed that the straight line $(\mathbf{x}_K, \mathbf{x}_L)$ is orthogonal to σ . For all $\sigma \in \mathcal{E}_{\text{ext}}$, there exists one unique cell K such that $\sigma \in \mathcal{E}_K$. Then we denote by \mathbf{x}_σ the projection of \mathbf{x}_K over the hyperplane containing σ , and we assume that \mathbf{x}_σ belongs to σ .*

In what follows, we denote by m_K the d -dimensional Lebesgue measure of the control volume $K \in \mathcal{T}$, and by m_σ the $(d-1)$ -Lebesgue measure of the edge $\sigma \in \mathcal{E}$. For all $\sigma \in \mathcal{E}_K$, we denote by $d_{K, \sigma} = d(\mathbf{x}_K, \mathbf{x}_\sigma)$. Since $\sigma = K|L$ is supposed to be orthogonal to $\mathbf{x}_K - \mathbf{x}_L$, then $d(\mathbf{x}_K, \mathbf{x}_L) = d_{K, \sigma} + d_{L, \sigma} =: d_\sigma$. We define the *transmissibilities* $(A_\sigma)_{\sigma \in \mathcal{E}}$ by

$$A_\sigma = \begin{cases} \frac{m_\sigma}{d_\sigma} & \text{if } \sigma = K|L \in \mathcal{E}_{\text{int}}, \\ \frac{m_\sigma}{d_{K, \sigma}} & \text{if } \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}. \end{cases}$$

The space of the degrees of freedom (including those prescribed by the boundary condition) is

$$\mathbb{X}_{\mathcal{T}} = \left\{ \mathbf{v} = (v_K, v_{\sigma})_{K \in \mathcal{T}, \sigma \in \mathcal{E}_{\text{ext}}} \right\} \simeq \mathbb{R}^{\#\mathcal{T} + \#\mathcal{E}_{\text{ext}}},$$

while the interior degrees of freedom (for which a nonlinear system has to be solved) are described by the space

$$\mathbb{X}_{\mathcal{T}, \text{int}} = \left\{ \mathbf{v} = (v_K)_{K \in \mathcal{T}} \right\} \simeq \mathbb{R}^{\#\mathcal{T}}.$$

2.1.2. Time and space-time discretizations.

Definition 2.2 (Time discretizations of $(0, T)$). *A time discretization of $(0, T)$ is given by an integer value N and a sequence of real values $0 = t^0 < t^1 < \dots < t^N = T$. For all $n \in \{1, \dots, N\}$ the time step is defined by $\Delta t^n = t^n - t^{n-1}$.*

Definition 2.3 (Space-time discretizations of Q). *A space-time discretization \mathcal{D} of Q is a family*

$$\mathcal{D} = (\mathcal{T}, \mathcal{E}, (\mathbf{x}_K)_{K \in \mathcal{T}}, (t^n)_{n \in \{0, \dots, N\}}),$$

where $(\mathcal{T}, \mathcal{E}, (\mathbf{x}_K)_{K \in \mathcal{T}})$ is an admissible mesh of Ω in the sense of Definition 2.1 and $(N, (t^n)_{n \in \{0, \dots, N\}})$ is a discretization of $(0, T)$ in the sense of Definition 2.2.

The spaces of the degrees of freedom are defined by

$$(21) \quad \mathbb{X}_{\mathcal{D}} = \left\{ \mathbf{v} = (v_K^n, v_{\sigma}^n)_{K \in \mathcal{T}, \sigma \in \mathcal{E}_{\text{ext}}, 1 \leq n \leq N} \right\} \simeq \mathbb{R}^{(\#\mathcal{T} + \#\mathcal{E}) \times N}.$$

and

$$(22) \quad \mathbb{X}_{\mathcal{D}, \text{int}} = \left\{ \mathbf{v} = (v_K^n)_{K \in \mathcal{T}, 1 \leq n \leq N} \right\} \simeq \mathbb{R}^{\#\mathcal{T} \times N}.$$

Let $\mathbf{v} = (v_K^n)_{K, n} \in \mathbb{X}_{\mathcal{D}}$, then we denote by $\mathbf{v}^n = (v_K^n)_K \in \mathbb{X}_{\mathcal{T}}$ for $n \in \{1, \dots, N\}$.

2.1.3. Reconstruction operators. Following the approach proposed in [20], we introduce reconstruction operators. First, we define the linear operator $\pi_{\mathcal{T}} : \mathbb{X}_{\mathcal{T}} \rightarrow L^{\infty}(\Omega)$ by

$$\pi_{\mathcal{T}} \mathbf{v}(\mathbf{x}) = v_K \quad \text{if } \mathbf{x} \in K, \quad \forall \mathbf{v} = (v_K, v_{\sigma})_{K \in \mathcal{T}, \sigma \in \mathcal{E}_{\text{ext}}}.$$

It is extended into the time-and-space reconstruction linear operator $\pi_{\mathcal{D}} : (\mathbb{X}_{\mathcal{T}})^N \rightarrow L^{\infty}(Q)$ by setting

$$\pi_{\mathcal{D}} \mathbf{v}(\mathbf{x}, t) = v_K^n \quad \text{if } (\mathbf{x}, t) \in K \times (t^{n-1}, t^n], \quad \forall \mathbf{v} = (v_K^n, v_{\sigma}^n)_{K \in \mathcal{T}, \sigma \in \mathcal{E}_{\text{ext}}, 1 \leq n \leq N} \in (\mathbb{X}_{\mathcal{T}})^N.$$

The study to be performed also requires the introduction of a so-called discrete gradient. We will remain sloppy about the construction of the discrete gradient. We only highlight the properties we will use in the sequel.

Lemma 2.4. *Let \mathcal{T} be an admissible discretization of Ω in the sense of Definition 2.1. There exists a linear operator $\nabla_{\mathcal{T}} : \mathbb{X}_{\mathcal{T}} \rightarrow L^{\infty}(\Omega)^d$ such that for all $\mathbf{v} = (v_K, v_{\sigma})_{K, \sigma}$ and $\mathbf{w} = (w_K, w_{\sigma})_{K, \sigma}$ in $\mathbb{X}_{\mathcal{T}}$, one has*

$$(23) \quad \int_{\Omega} \nabla_{\mathcal{T}} \mathbf{v} \cdot \nabla_{\mathcal{T}} \mathbf{w} d\mathbf{x} = \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} A_{\sigma} (v_K - v_L) (w_K - w_L) + \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K, \text{ext}}} A_{\sigma} (v_K - v_{\sigma}) (w_K - w_{\sigma}).$$

Moreover, let $(\mathcal{T}_m)_{m \geq 1}$ be a sequence of admissible discretizations of Ω in the sense of Definition 2.1 such that $\text{size}(\mathcal{T}_m)$ tends to 0 while $\text{reg}(\mathcal{T}_m)$ remains bounded as m tends to ∞ , and let $(\mathbf{v}_m)_{m \geq 1}$ be a family such that $\mathbf{v}_m \in \mathbb{X}_{\mathcal{T}_m}$ for all $m \geq 1$ and

$$\|\pi_{\mathcal{T}} \mathbf{v}_m\|_{L^2(\Omega)} + \|\nabla_{\mathcal{T}_m} \mathbf{v}_m\|_{L^2(\Omega)^d} \leq C, \quad \forall m \geq 1,$$

then there exists $v \in H^1(\Omega)$ such that, up to an unlabeled subsequence, one has

$$(24) \quad \pi_{\mathcal{T}_m} \mathbf{v}_m \xrightarrow{m \rightarrow \infty} v \quad \text{in } L^2(\Omega),$$

and

$$(25) \quad \nabla_{\mathcal{T}_m} \mathbf{v}_m \xrightarrow{m \rightarrow \infty} \nabla v \quad \text{weakly in } L^2(\Omega)^d.$$

Additionally, if $\varphi \in C_c^\infty(\Omega)$ is discretized into $\varphi_m = (\varphi_K)_{K \in \mathcal{T}_m} \in \mathbb{X}_{\mathcal{T}_m}$ by setting

$$\varphi_K = \frac{1}{m_K} \int_K \varphi(\mathbf{x}) d\mathbf{x}, \quad \forall K \in \mathcal{T}_m,$$

then $\nabla_{\mathcal{T}_m} \varphi_m$ converges strongly in $L^2(\Omega)^d$ towards $\nabla \varphi$ as m tends to $+\infty$.

Note that the discrete gradient reconstruction used in [2, §4.1] (which is the usual one for two-point flux approximations of diffusion operators) does not meet the requirements of Lemma 2.4 since (23) is not fulfilled. However, a reconstruction as prescribed by Lemma 2.4 can be obtained as a particular case of the so-called SUSHI scheme on “super-admissible meshes” (cf. [24, Lemma 2.1]).

Finally, the reconstruction operators $\pi_{\mathcal{T}} : \mathbb{X}_{\mathcal{T}} \rightarrow L^\infty(\Omega)$ and $\nabla_{\mathcal{T}} : \mathbb{X}_{\mathcal{T}} \rightarrow L^\infty(\Omega)^d$ are extended to the space-times framework into $\pi_{\mathcal{D}} : \mathbb{X}_{\mathcal{D}} \rightarrow L^\infty(Q)$ and $\nabla_{\mathcal{D}} : \mathbb{X}_{\mathcal{D}} \rightarrow L^\infty(Q)^d$ defined for all $\mathbf{v} = (\mathbf{v}^n)_{1 \leq n \leq N} \in \mathbb{X}_{\mathcal{D}}$ by

$$(26) \quad \pi_{\mathcal{D}} \mathbf{v}(\cdot, t) = \pi_{\mathcal{T}} \mathbf{v}^n, \quad \nabla_{\mathcal{D}} \mathbf{v}(\cdot, t) = \nabla_{\mathcal{T}} \mathbf{v}^n, \quad \forall t \in (t^{n-1}, t^n].$$

2.2. The implicit finite volume scheme. The initial data s_0 is discretized into $\mathbf{s}^0 = (s_K^0)_{K \in \mathcal{T}} \in \mathbb{X}_{\mathcal{T}, \text{int}}$ by setting

$$(27) \quad s_K^0 = \frac{1}{m_K} \int_K s_0(\mathbf{x}) d\mathbf{x}, \quad \forall K \in \mathcal{T}.$$

Notice that $0 \leq s_K^0 \leq 1$ since $0 \leq s_0 \leq 1$. We define $\boldsymbol{\tau}^0 = (\tau_K^0)_K \in \mathbb{X}_{\mathcal{T}, \text{int}}$ as

$$(28) \quad \boldsymbol{\tau}^0 = \mathbf{s}^{-1}(\mathbf{s}^0),$$

so that

$$s(\tau_K^0) = s_K^0, \quad \forall K \in \mathcal{T}.$$

The boundary condition $\tau_D \in C^1(\overline{Q})$ is discretized by setting for all $n \in \{0, \dots, N\}$

$$(29) \quad \tau_{D, \sigma}^n = \tau_D(\mathbf{x}_\sigma, t^n), \quad \forall \sigma \in \mathcal{E}_{\text{ext}}, \quad \text{and} \quad \tau_{D, K}^n = \tau_D(\mathbf{x}_K, t^n), \quad \forall K \in \mathcal{T}.$$

Then we denote by $\boldsymbol{\tau}_D^n = (\tau_{D, K}^n, \tau_{D, \sigma}^n)_{K \in \mathcal{T}, \sigma \in \mathcal{E}_{\text{ext}}} \in \mathbb{X}_{\mathcal{T}}$, and by $\boldsymbol{\tau}_D = (\boldsymbol{\tau}_D^n)_{1 \leq n \leq N} \in \mathbb{X}_{\mathcal{D}}$. It follows from Formula (23) and from the regularity of τ_D that

$$(30) \quad \begin{aligned} \int_{\Omega} |\nabla_{\mathcal{T}} \boldsymbol{\tau}_D^n|^2 d\mathbf{x} &= \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} A_\sigma (\tau_{D, K}^n - \tau_{D, L}^n)^2 + \sum_{\sigma \in \mathcal{E}_{\text{ext}}} A_\sigma (\tau_{D, K}^n - \tau_{D, \sigma}^n)^2 \\ &\leq \|\nabla \tau_D\|_\infty^2 \left(\sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} m_\sigma d_\sigma + \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m_\sigma d_{K, \sigma} \right) = dm_\Omega \|\nabla \tau_D\|_\infty^2. \end{aligned}$$

Let $n \geq 1$. Assume that the state $\boldsymbol{\tau}^{n-1} = (\tau_K^{n-1})_K \in \mathbb{X}_{\mathcal{T}}$ is known. The implicit finite volume scheme is obtained by writing the local conservation of the volume of each fluid on the control volumes, i.e.,

$$(31) \quad \frac{s(\tau_K^n) - s(\tau_K^{n-1})}{\Delta t^n} m_K + \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^n = 0, \quad \forall K \in \mathcal{T},$$

where $F_{K,\sigma}^n$ denotes the outward w.r.t. K flux across the edge σ at time step t^n . Denote by $\mathbf{n}_{K,\sigma}$ the outward w.r.t. K normal to σ , and by $g_{K,\sigma} = \mathbf{g} \cdot \mathbf{n}_{K,\sigma}$ for all $\sigma \in \mathcal{E}_K$ and all $K \in \mathcal{T}$. Denote by

$$\tau_{K,\sigma}^n = \begin{cases} \tau_L^n & \text{if } \sigma = K|L \in \mathcal{E}_{\text{int}}, \\ \tau_{D,\sigma}^n & \text{if } \sigma \in \mathcal{E}_{K,\text{ext}}, \end{cases}$$

then the fluxes $F_{K,\sigma}^n$ across $\sigma \in \mathcal{E}_K$ is defined by

$$(32) \quad F_{K,\sigma}^n = m_\sigma \left(\lambda(s(\tau_K^n)) g_{K,\sigma}^+ - \lambda(s(\tau_{K,\sigma}^n)) g_{K,\sigma}^- \right) + A_\sigma (u(\tau_K^n) - u(\tau_{K,\sigma}^n)).$$

Note in particular that the scheme is locally conservative, i.e.,

$$F_{K,\sigma}^n + F_{L,\sigma}^n = 0, \quad \forall \sigma = K|L \in \mathcal{E}_{\text{int}}.$$

Combining (31)–(32), the scheme can be rewritten in a condensed form as

$$(33) \quad \mathcal{R}_K \left(\tau_K^n, \tau_K^{n-1}, (\tau_L^n)_{L \neq K}, (\tau_{D,\sigma}^n)_{\sigma \in \mathcal{E}_{K,\text{ext}}} \right) = 0, \quad \forall K \in \mathcal{T},$$

where \mathcal{R}_K is nondecreasing w.r.t. its first argument and nonincreasing w.r.t. the others thanks to the monotonicity of the functions λ, s and u .

It is worth noticing that

$$(34) \quad \sum_{\sigma \in \mathcal{E}_K} m_\sigma g_{K,\sigma} = 0, \quad \forall K \in \mathcal{T}.$$

Therefore, the convective flux balance can be reformulated, yielding

$$(35) \quad \begin{aligned} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^n &= \sum_{\sigma \in \mathcal{E}_{K,L}} m_\sigma g_{K,\sigma}^- (\lambda(s(\tau_K^n)) - \lambda(s(\tau_{K,\sigma}^n))) \\ &\quad + \sum_{\sigma=K|L} A_\sigma (u(\tau_K^n) - u(\tau_{K,\sigma}^n)), \quad \forall K \in \mathcal{T}. \end{aligned}$$

2.3. Existence and uniqueness of the solution to the scheme. In this section, we analyze the system (33) obtained for a fixed admissible discretization \mathcal{D} of Q . In what follows, we denote by

$$a \top b = \max(a, b) \quad \text{and} \quad a \perp b = \min(a, b), \quad \forall (a, b) \in \mathbb{R}^2.$$

The following Lemma is a discrete counterpart of the L^1 -contraction principle (19) on the exact solution.

Lemma 2.5. *Let $\boldsymbol{\tau}^{n-1}$ and $\widehat{\boldsymbol{\tau}}^{n-1}$ be two elements of $\mathbb{X}_{\mathcal{T},\text{int}}$, and let $\boldsymbol{\tau}^n$ and $\widehat{\boldsymbol{\tau}}^n$ in $\mathbb{X}_{\mathcal{T},\text{int}}$ be two corresponding solutions, then*

$$(36) \quad \int_{\Omega} |\pi_{\mathcal{T}} s(\boldsymbol{\tau}^n) - \pi_{\mathcal{T}} s(\widehat{\boldsymbol{\tau}}^n)| \, d\mathbf{x} \leq \int_{\Omega} |\pi_{\mathcal{T}} s(\boldsymbol{\tau}^{n-1}) - \pi_{\mathcal{T}} s(\widehat{\boldsymbol{\tau}}^{n-1})| \, d\mathbf{x}.$$

Proof. It follows from the monotonicity of \mathcal{R}_K that

$$\begin{aligned}\mathcal{R}_K \left(\tau_K^n, \tau_K^{n-1} \top \hat{\tau}_K^{n-1}, (\tau_L^n \top \hat{\tau}_L^n)_{L \neq K}, (\tau_{D,\sigma})_{\sigma \in \mathcal{E}_{K,\text{ext}}} \right) &\leq 0, \\ \mathcal{R}_K \left(\hat{\tau}_K^n, \tau_K^{n-1} \top \hat{\tau}_K^{n-1}, (\tau_L^n \top \hat{\tau}_L^n)_{L \neq K}, (\tau_{D,\sigma})_{\sigma \in \mathcal{E}_{K,\text{ext}}} \right) &\leq 0.\end{aligned}$$

Since $\tau_K^n \top \hat{\tau}_K^n$ is either equal to τ_K^n or to $\hat{\tau}_K^n$, one has

$$(37) \quad \mathcal{R}_K \left(\tau_K^n \top \hat{\tau}_K^n, \tau_K^{n-1} \top \hat{\tau}_K^{n-1}, (\tau_L^n \top \hat{\tau}_L^n)_{L \neq K}, (\tau_{D,\sigma})_{\sigma \in \mathcal{E}_{K,\text{ext}}} \right) \leq 0.$$

Similar calculations lead to

$$(38) \quad \mathcal{R}_K \left(\tau_K^n \perp \hat{\tau}_K^n, \tau_K^{n-1} \perp \hat{\tau}_K^{n-1}, (\tau_L^n \perp \hat{\tau}_L^n)_{L \neq K}, (\tau_{D,\sigma})_{\sigma \in \mathcal{E}_{K,\text{ext}}} \right) \geq 0.$$

Summing (37) with (38) and over $K \in \mathcal{T}$ yields (36). \square

Lemma 2.6. *Given $\tau^{n-1} \in \mathbb{X}_{\mathcal{T}}$, then there exists at most one solution $\tau^n \in \mathbb{X}_{\mathcal{T}}$ to the scheme (31)–(32).*

Proof. As a direct consequence of Lemma 2.5, $s(\tau_K^n) = s(\hat{\tau}_K^n)$ for all $K \in \mathcal{T}$. Subtracting the system yielding $\hat{\tau}^n$ to the one corresponding to τ^n leads to

$$\sum_{\sigma=K|L \in \mathcal{E}_K} A_{\sigma} (w_K^n - w_L^n) + \sum_{\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K} A_{\sigma} w_K^n = 0, \quad \forall K \in \mathcal{T},$$

where we have set $w_K^n = u(\tau_K^n) - u(\hat{\tau}_K^n)$. It follows from classical arguments that $w^n = (w_K^n)_K = \mathbf{0}_{\mathbb{X}_{\mathcal{T}}}$. Bearing the nondegeneracy condition (11) in mind, we get that $\tau^n = \hat{\tau}^n$. \square

The following lemma shows that the solution τ^n to the scheme is always greater than τ^* . Therefore, the extension (10) we chose for the function u does not affect the result.

Lemma 2.7. *Let $\tau^n \in \mathbb{X}_{\mathcal{T}}$ be the solution to (31)–(32), then $\tau_K^n \geq \tau_*$ for all $K \in \mathcal{T}$.*

Proof. There is nothing to prove if $\tau_* = -\infty$, hence let us assume that τ_* is finite. Let K be a cell such that $\tau_K^n \leq \tau_{K,\sigma}^n$ for all $\sigma \in \mathcal{E}_K$, and assume that $\tau_K^n < \tau_*$. Since $s(\tau_K^n) = 0$ and $\lambda(s(\tau_K^n)) = 0$, one gets that

$$\sum_{\sigma \in \mathcal{E}_K} A_{\sigma} (u(\tau_K^n) - u(\tau_{K,\sigma}^n)) \geq 0.$$

The extension (10) of u ensures that $u(\tau_K^n) < u(\tau_*) \leq u(\tau_{D,\sigma}^n)$. The left-hand side of the above relation is therefore negative, hence a contradiction with the assumption $\tau_K^n < \tau_*$. \square

Let $n \in \{0, \dots, N\}$, then define

$$e_K^n(\tau) = \int_{\tau_{D,K}^n}^{\tau} (a - \tau_{D,K}^n) s'(a) da = \int_{\tau_{D,K}^n}^{\tau} (s(\tau) - s(a)) da \geq 0, \quad \forall \tau \in \mathbb{R},$$

and $\mathfrak{E}^n : \mathbb{X}_{\mathcal{T}} \rightarrow \mathbb{R}_+$ by

$$\mathfrak{E}^n(\tau) = \sum_{K \in \mathcal{T}} e_K^n(\tau_K) m_K, \quad \forall \tau = (\tau_K)_{K \in \mathcal{T}}.$$

It is easy to verify (see e.g. [11]) that

$$(39) \quad 0 \leq \mathfrak{E}^n(\tau) \leq m_{\Omega} (\|1 - s\|_{L^1(\mathbb{R}_+)} + \|s\|_{L^1(\mathbb{R}_-)}), \quad \forall \tau \in \mathbb{X}_{\mathcal{T}}.$$

Moreover, it follows from the C^1 regularity of τ_D and from the fact that $0 \leq s \leq 1$ that

$$(40) \quad |\mathfrak{E}^n(\tau) - \mathfrak{E}^{n-1}(\tau)| \leq \Delta t^n m_\Omega \|\partial_t \tau_D\|_\infty.$$

We define the Lipschitz continuous function $\xi : \mathbb{R}_+ \rightarrow \mathbb{R}$ by

$$(41) \quad \xi(\tau) = \int_0^\tau \sqrt{\lambda(s(a))p'(a)} da = \int_0^\tau \sqrt{u'(a)} da, \quad \forall \tau \in \mathbb{R},$$

then it follows from Cauchy-Schwartz inequality that

$$(42) \quad (a - b)(u(a) - u(b)) \geq (\xi(a) - \xi(b))^2, \quad \forall (a, b) \in \mathbb{R}^2.$$

Moreover, the Lipschitz continuity of ξ implies that

$$u(\tau) \leq \|\xi'\|_\infty \xi(\tau), \quad \forall \tau \in \mathbb{R}_+,$$

hence it follows from Assumption (13) that

$$(43) \quad \tau \leq C (\xi(\tau) + 1), \quad \forall \tau \in \mathbb{R}_+,$$

then, in particular, one has

$$(44) \quad \lim_{\tau \rightarrow \infty} \xi(\tau) = +\infty.$$

Lemma 2.8. *Let τ^n be a solution to the scheme (31)–(32). Then there exists C_1 depending only on Ω and τ_D such that the following estimate holds:*

$$(45) \quad \mathfrak{E}^n(\tau^n) + \Delta t^n \left(C_1 + \frac{1}{2} \int_\Omega |\nabla_{\mathcal{T}} \xi(\tau)|^2 d\mathbf{x} \right) \leq \mathfrak{E}^{n-1}(\tau^{n-1}).$$

Proof. We multiply the equation (31) by $\Delta t^n (\tau_K^n - \tau_{D,K}^n)$ and sum over $K \in \mathcal{T}$. Using (35), this provides

$$(46) \quad T_1 + \Delta t^n (T_2 + T_3) = 0,$$

where

$$\begin{aligned} T_1 &= \sum_{K \in \mathcal{T}} (s(\tau_K^n) - s(\tau_K^{n-1})) (\tau_K^n - \tau_{D,K}^n) m_K, \\ T_2 &= \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m_\sigma g_{K,\sigma}^- (\lambda(s(\tau_K^n)) - \lambda(s(\tau_{K,\sigma}^n))) (\tau_K^n - \tau_{D,K}^n), \\ T_3 &= \int_\Omega \nabla_{\mathcal{T}} u(\tau^n) \cdot \nabla_{\mathcal{T}} (\tau^n - \tau_D^n) d\mathbf{x}. \end{aligned}$$

It follows from the convexity of $e \circ s^{-1}$ (see e.g. [9, Proposition 3.7]) that

$$T_1 \geq \mathfrak{E}^n(\tau^n) - \mathfrak{E}^n(\tau^{n-1}).$$

Then thanks to (40), one gets that

$$(47) \quad T_1 \geq \mathfrak{E}^n(\tau^n) - \mathfrak{E}^{n-1}(\tau^{n-1}) - \Delta t^n m_\Omega \|\partial_t \tau_D\|_\infty.$$

The term T_2 can be estimated following the path of [21]. Denote by $\Phi : \mathbb{R} \rightarrow \mathbb{R}_+$ the function defined by

$$\Phi(\tau) = \int_0^\tau a s'(a) \lambda'(s(a)) da, \quad \forall \tau \in \mathbb{R}.$$

Since $s(0) = 1$, one has $s'(a) = 0$ for all $a \geq 0$, and thus

$$\Phi(\tau) = 0 \quad \text{if } \tau > 0, \quad \text{and} \quad \Phi'(\tau) \leq 0, \quad \forall \tau \leq 0.$$

Moreover, since $s \in L^1(\mathbb{R}_-)$, one has

$$|\tau|\lambda(s(\tau)) \leq \|\lambda'\|_\infty |\tau|s(\tau) \xrightarrow{\tau \searrow \tau_*} 0.$$

Therefore, for all $\tau \leq 0$, one has

$$\Phi(\tau) = \tau\lambda(s(\tau)) + \int_\tau^0 \lambda(s(a))da \leq \int_\tau^0 \lambda(s(a))da + C.$$

Thanks to (9) and to Assumption (14), one has $\int_\tau^0 \lambda(s(a))da \leq C$, hence Φ is bounded. Simple calculations show that for all $(a, b) \in \mathbb{R}^2$, one has

$$b(\lambda(s(b)) - \lambda(s(a))) = \Phi(b) - \Phi(a) + \int_a^b (\lambda(s(r)) - \lambda(s(a)))dr \geq \Phi(b) - \Phi(a).$$

Rewriting

$$(48) \quad T_2 = T_{21} + T_{22}$$

with

$$\begin{aligned} T_{21} &= \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m_\sigma g_{K,\sigma}^- (\lambda(s(\tau_K^n)) - \lambda(s(\tau_{K,\sigma}^n))) \tau_K^n, \\ T_{22} &= - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m_\sigma g_{K,\sigma}^- (\lambda(s(\tau_K^n)) - \lambda(s(\tau_{K,\sigma}^n))) \tau_{D,K}^n, \end{aligned}$$

we get that

$$\begin{aligned} T_{21} &\geq \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m_\sigma g_{K,\sigma}^- (\Phi(s(\tau_K^n)) - \Phi(s(\tau_{D,K}^n))) \\ &\geq \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K,\text{ext}}} m_\sigma (g_{K,\sigma}^+ \Phi(s(\tau_K^n)) - g_{K,\sigma}^- \Phi(s(\tau_{D,K}^n))). \end{aligned}$$

Using the boundedness of Φ , we get that

$$(49) \quad T_{21} \geq -m_{\partial\Omega} |\mathbf{g}| \|\Phi\|_\infty.$$

On the other hand, a classical reorganization of the term T_{22} provides

$$\begin{aligned} T_{22} &= \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} m_\sigma (g_{K,\sigma}^+ \lambda(s(\tau_K^n)) - g_{L,\sigma}^+ \lambda(s(\tau_L^n))) (\tau_{D,K}^n - \tau_{D,L}^n) \\ &\quad - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{K,\text{ext}}} m_\sigma (g_{K,\sigma}^+ \lambda(s(\tau_K^n)) - g_{K,\sigma}^- \lambda(s(\tau_{D,\sigma}^n))) \tau_{D,K}^n. \end{aligned}$$

Therefore, it follows from the regularity of τ_D and from the boundedness of λ that

$$(50) \quad T_{22} \geq -|\mathbf{g}| \|\lambda\|_\infty (dm_\Omega \|\nabla \tau_D\|_\infty - m_{\partial\Omega} \|\tau_D\|_\infty).$$

Finally, it results from (42), from the relation $u' = (\xi')^2$, and from (30) that

$$\begin{aligned} T_3 &\geq \|\nabla \tau \xi(\tau^n)\|_{L^2(\Omega)^d}^2 - \|\nabla \tau u(\tau^n)\|_{L^2(\Omega)^d} \|\nabla \tau \tau_D^n\|_{L^2(\Omega)^d} \\ &\geq \|\nabla \tau \xi(\tau^n)\|_{L^2(\Omega)^d}^2 - \sqrt{\|u'\|_\infty} \|\nabla \tau \xi(\tau^n)\|_{L^2(\Omega)^d} \|\nabla \tau \tau_D^n\|_{L^2(\Omega)^d}. \end{aligned}$$

Therefore, it follows from (30) that

$$(51) \quad T_3 \geq \frac{1}{2} \int_\Omega |\nabla \tau \xi(\tau^n)|^2 d\mathbf{x} - \frac{\|u'\|_\infty dm_\Omega \|\nabla \tau_D\|_\infty}{2}.$$

Putting (47)–(51) in (46) ends the proof of Lemma 2.8. \square

Proposition 2.9. *Let $\tau^{n-1} \in \mathbb{X}_{\mathcal{T},\text{int}}$, there exists a unique solution $\tau^n \in \mathbb{X}_{\mathcal{T},\text{int}}$ to the scheme (31)–(32). Moreover, it satisfies $\tau_K^n \geq \tau_*$ for all $K \in \mathcal{T}$.*

Proof. The uniqueness of the solution was proven at Lemma 2.6 while the fact that $\tau^n \geq \tau_*$ was the purpose of Lemma 2.7. Therefore, it only remains to show the existence of a solution. It follows from Estimate (45) that

$$\|\nabla_{\mathcal{T}}(\xi(\tau^n) - \xi(\tau_D^n))\|_{L^2(\Omega)^d}^2 \leq 2 \frac{\mathfrak{E}^{n-1}(\tau^{n-1})}{\Delta t^n} + 2C_1 + 4\|u'\|_{\infty} \|\nabla \tau_D\|_{\infty}^2.$$

Since ξ is coercive (44), there exists C depending on the data (among which τ^{n-1} , τ_D , the mesh \mathcal{T} and the time step Δt^n) such that

$$|\tau_K^n| \leq C, \quad \forall K \in \mathcal{T}.$$

This estimate allows us to make use of a topological degree argument (see [34, 16, 21]) to prove the existence of one solution to the scheme. \square

The convergence of the scheme can be proved following the path proposed in [25]. Enhanced convergence properties can be obtained thanks to the recent contribution [19]. But this is not the goal of this paper. We are mainly interested in the practical computation of the approximate solution at fixed discretization parameters. In particular, we focus on the behavior the Newton's method.

3. ABOUT THE NEWTON METHOD

The numerical scheme (31)–(32) amounts for all $n \in \{1, \dots, N\}$ to the nonlinear system

$$(52) \quad \mathcal{F}_n(\tau^n) = (f_K^n(\tau^n))_{K \in \mathcal{T}} = \mathbf{0}, \quad \text{with } \mathcal{F}_n \in \mathcal{C}^2(\mathbb{R}^{\#\mathcal{T}}; \mathbb{R}^{\#\mathcal{T}}),$$

where

$$\begin{aligned} f_K(\tau) &= (s(\tau_K) - s_K^{n-1}) \\ &+ \frac{\Delta t^n}{m_K} \sum_{\sigma \in \mathcal{E}_K} \left(m_{\sigma} g_{K,\sigma}(\lambda(s(\tau_K)) - \lambda(s(\tau_{K,\sigma}^n))) + A_{\sigma} (u(\tau_K) - u(\tau_{K,\sigma}^n)) \right). \end{aligned}$$

Assume that the Jacobian matrix $\mathbb{J}_{\mathcal{F}_n}(\tau^n)$ of \mathcal{F}_n at τ^n is not singular (this will be shown for nondegenerate parametrizations (11), cf. Proposition 3.7). Approximating the solutions to the system (20) with the Newton method consists in the construction of a sequence $(\tau^{n,k})_{k \geq 0}$ defined by

$$(53) \quad \begin{cases} \tau^{n,0} = \tau^{n-1}, \\ \tau^{n,k+1} = \tau^{n,k} - [\mathbb{J}_{\mathcal{F}_n}(\tau^{n,k})]^{-1} \mathcal{F}_n(\tau^{n,k}). \end{cases}$$

If the method converges, the solution τ^n is then defined as

$$(54) \quad \tau^n = \lim_{k \rightarrow \infty} \tau^{n,k}.$$

Moreover, the convergence speed is asymptotically quadratic if $\mathcal{F}_n \in \mathcal{C}^2$, i.e.,

$$(55) \quad \|\tau^{n,k} - \tau^n\| \leq C \|\tau^{n,k-1} - \tau^n\|^2, \quad \forall k \geq k_{\star} \text{ large enough}$$

where, denoting by $\mathcal{V}(\tau^n)$ a neighborhood of τ^n in $\mathbb{X}_{\mathcal{T}}$, the quantity C (as well as k_{\star}) depends on

$$(56) \quad \sup_{\tau \in \mathcal{V}(\tau^n)} \left\| [\mathbb{J}_{\mathcal{F}_n}(\tau)]^{-1} \right\| \quad \text{and} \quad \sup_{\tau \in \mathcal{V}(\tau^n)} \|D^2 \mathcal{F}_n(\tau)\|.$$

Since τ^n is unknown, a sufficient condition to ensure that (55) holds for some $C > 0$ is the following uniform non-degeneracy condition

$$(57) \quad \sup_{\tau \in \mathbb{R} \setminus \tau} \left\| [\mathbb{J}_{\mathcal{F}_n}(\tau)]^{-1} \right\| < \infty.$$

We cite here a simplified version of the so-called Newton-Kantorovich theorem (see, e.g., [32, 37, 38]). We refer to [29] for a quantitative version of the theorem, and to [41] to a non-smooth version.

Theorem 3.1 (Newton-Kantorovich theorem). *Assume that there exists two positive quantities C_2 and C_3 such that*

$$(58) \quad \sup_{\tau \in \mathbb{X}_{\mathcal{T}}} \|\mathbb{J}_{\mathcal{F}_n}(\tau)\| \leq C_2 \quad \text{and} \quad \sup_{\tau \in \mathbb{X}_{\mathcal{T}}} \left\| [\mathbb{J}_{\mathcal{F}_n}(\tau)]^{-1} \right\| \leq C_3,$$

then there exists $\rho > 0$ such that

$$\|\tau^{n,0} - \tau^n\| \leq \rho \quad \implies \quad \tau^{n,k} \xrightarrow[k \rightarrow \infty]{} \tau^n.$$

Our strategy in the sequel is to prove that a non-degenerate parametrization (in the sense of (11)) yields estimates (58) in the subordinate matrix 1-norm.

Remark 3.2. *The assumption $\mathcal{F} \in \mathcal{C}^1$ can be relaxed. More precisely, Newton method can be extended to the case where \mathcal{F} is merely semi-smooth [15] following the way proposed by Qi and Sun in [41]. Quadratic convergence is preserved in this nonsmooth case provided \mathcal{F}_n is semi-smooth of order 1, cf. [41, Theorem 3.2]. Our study can be extended to this more general case, but, for the sake of simplicity, we have chosen to reduce our presentation to the classical smooth framework.*

3.1. Some technical lemmas related to M -matrices. Because of the elliptic degeneracy of the problem when $\tau \geq 0$, we cannot apply the results of [26] to get estimates on the Jacobian matrix $\mathbb{J}_{\mathcal{F}_n}$. We need to introduce some technical material to circumvent this difficulty. Let us first define that notion of δ -transmissive path (see also [9, 10])

Definition 3.3 (δ -transmissive path). *Let $\delta > 0$, let $\mathbb{A} = (a_{ij})_{1 \leq i, j \leq N} \in \mathcal{M}_N(\mathbb{R})$ and let $(i, j) \in \{1, \dots, N\}$.*

- (1) *A row-wise δ -transmissive path $\mathcal{P}(i, j)$ of length L from i to j associated to the matrix \mathbb{A} consists in a list $\{k_0, \dots, k_L\}$ with*
 - (i) $k_0 = i$, $k_L = j$, and $k_p \neq k_q$ if $p \neq q$;
 - (ii) *for all $p \in \{0, \dots, L-1\}$, one has $a_{k_p k_{p+1}} < -\delta$.*
- (2) *The list $\{k_0, \dots, k_L\}$ is a column-wise δ -transmissive path associated to the matrix \mathbb{A} if it is a row-wise δ -transmissive path associated to \mathbb{A}^T .*

Definition 3.4 ((δ, Δ) - M -matrices). *Let $\delta, \Delta > 0$ be such that $\Delta > \delta$.*

- (1) *A matrix $\mathbb{A} = (a_{ij})_{1 \leq i, j \leq N} \in \mathcal{M}_N(\mathbb{R})$ is said to be a row-wise (δ, Δ) - M -matrix if*
 - (i) *for all $i \in \{1, \dots, N\}$, one has $\delta \leq a_{ii} \leq \Delta$, $a_{ij} \leq 0$ for all $j \neq i$ and $\sum_{j=1}^N a_{ij} \geq 0$;*
 - (ii) *the set $\mathcal{I}_{\delta}(\mathbb{A}) = \left\{ i \in \{1, \dots, N\} \mid \sum_{j=1}^N a_{ij} \geq \delta \right\}$ is not empty, and for all $i \in \mathcal{I}_{\delta}(\mathbb{A})^c = \{1, \dots, N\} \setminus \mathcal{I}_{\delta}(\mathbb{A})$, there exists a δ -transmissive path $\mathcal{P}(i, j)$ with $j \in \mathcal{I}_{\delta}(\mathbb{A})$.*

- (2) A matrix $\mathbb{A} = (a_{ij})_{1 \leq i, j \leq N} \in \mathcal{M}_N(\mathbb{R})$ is said to be a column-wise (δ, Δ) -M-matrix if \mathbb{A}^T is a row-wise (δ, Δ) -M-matrix.

It is well known that M-matrices are invertible matrices. The goal of Lemma 3.5 and of Corollary 3.6 is to get uniform estimates on the inverse of a uniform M-matrix. Let us stress that the estimates we obtain are far from being optimal in the applications we have in mind, namely the Finite Volume discretization of Richards' equation.

Lemma 3.5. *Let $\mathbb{A} \in \mathcal{M}_N(\mathbb{R})$ be a row-wise (δ, Δ) -M-matrix, then there exists C depending only on δ, Δ and N such that $\|\mathbb{A}^{-1}\|_\infty \leq C$.*

Proof. For the ease of reading, we denote $\|\cdot\|$ instead of $\|\cdot\|_\infty$, and \mathcal{I}_δ instead of $\mathcal{I}_\delta(\mathbb{A})$. The property $\|\mathbb{A}^{-1}\| = \frac{1}{\alpha}$ is equivalent to

$$(59) \quad \frac{1}{\|\mathbb{A}^{-1}\|} = \max_{\|z\|=1} \|\mathbb{A}z\| = \alpha.$$

Let $z = (z_i)_{1 \leq i \leq N} \in \mathbb{R}^N$ with $\|z\| = 1$. We assume, without loss of generality that there exists $i \in \{1, \dots, N\}$ such that $z_i = 1$. In the sequel, we denote by \mathcal{L} the maximum length of a transmissive path, i.e.,

$$(60) \quad \mathcal{L} := \max_{j \in \mathcal{I}_\delta(\mathbb{A})^c} \min_{j \in \mathcal{I}_\delta(\mathbb{A})} \text{length}(\mathcal{P}(i, j)) \leq N - 1.$$

Assume first that $i \in \mathcal{I}_\delta$, then $(\mathbb{A}z)_i = \sum_j a_{ij}z_j \geq \delta$, hence $\|\mathbb{A}^{-1}\| \leq \frac{1}{\delta}$. Assume now that $i \in \mathcal{I}_\delta^c$, and let $\{k_p, 0 \leq p \leq L\}$ be a transmissive path with $k_0 = i$, $k_L \in \mathcal{I}_\delta$ and $L \leq \mathcal{L}$. Denote by

$$c_p = \left(\left(\frac{\Delta}{\delta} \right)^p - 1 \right) \frac{1}{\Delta - \delta}.$$

Let us show by induction that

$$(61) \quad z_{k_p} \geq 1 - c_p \alpha \quad \forall p \in \{0, \dots, L\}.$$

Since $z_{k_0} = 1$, the relation (61) holds for $p = 0$. Now suppose that (61) holds for some $p \in \{0, \dots, N - 1\}$. One knows from (59) that $\sum_{j=1}^N a_{k_p j} z_j \leq \alpha$, hence, using (61), that $\sum_{j=1}^N a_{k_p j} \geq 0$, and that $a_{k_p j} z_j \geq a_{k_p j}$ for all $j \notin \{p, p+1\}$, one gets that

$$-a_{k_p k_{p+1}} (1 - z_{k_{p+1}}) \leq \alpha (1 + c_p a_{k_p k_p}).$$

Since $a_{k_p k_p} \leq \Delta$ and $a_{k_p k_{p+1}} \leq -\delta$, this leads to

$$z_{k_{p+1}} \geq 1 - \alpha \frac{1 + \Delta c_p}{\delta} = 1 - \alpha c_{p+1},$$

so that the proof of (61) is complete. As a consequence, one gets that $z_{k_L} \geq 1 - \alpha c_L$, and thus that

$$a_{k_L k_L} (1 - c_L \alpha) + \sum_{j \neq k_L} a_{k_L j} \leq \alpha.$$

Using that $k_L \in \mathcal{I}_\delta$, we obtain that $(\Delta c_L + 1) \alpha \geq \delta$. Therefore,

$$\alpha \geq \frac{1}{c_{L+1}} \geq \frac{1}{c_{\mathcal{L}+1}}$$

since the length L of the path $\{k_0, \dots, k_L\}$ is bounded by \mathcal{L} defined by (60). \square

Corollary 3.6. *Let $\mathbb{A} \in \mathcal{M}_N(\mathbb{R})$ be a column-wise (δ, Δ) -M-matrix, then there exists C depending only on δ, Δ and N such that $\|\mathbb{A}^{-1}\|_1 \leq C$.*

3.2. A uniform non-degeneracy result.

Proposition 3.7. *Assume that s and u satisfy the non-degeneracy condition (11), then there exist C_2 depending only on $\alpha_\star, \mathbf{g}, \|\lambda'\|_\infty, \Delta t^n$ and \mathcal{T} , and C_3 depending only on $\alpha_\star, \alpha^\star, \Delta t^n$ and \mathcal{T} such that*

$$\|\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})\|_1 \leq C_2, \quad \|\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})\|_1^{-1} \leq C_3, \quad \forall \boldsymbol{\tau} \in \mathbb{R}^{\#\mathcal{T}}.$$

Proof. Let us first make the Jacobian matrix $\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau}) = (j_{KL}^n(\boldsymbol{\tau}))_{K,L \in \mathcal{T}}$ explicit:

$$\begin{aligned} j_{KK}^n(\boldsymbol{\tau}) &= s'(\tau_K) + \frac{\Delta t^n}{m_K} \sum_{\sigma \in \mathcal{E}_K} \left(m_\sigma g_{K,\sigma}^+ \lambda'(s(\tau_K)) s'(\tau_K) + A_\sigma u'(\tau_K) \right), \\ j_{LK}^n(\boldsymbol{\tau}) &= -\frac{\Delta t^n}{m_K} \left(m_\sigma g_{K,\sigma}^+ \lambda'(s(\tau_K)) s'(\tau_K) + A_\sigma u'(\tau_K) \right) \quad \text{where } \sigma = K|L. \end{aligned}$$

Proving that $\|\mathbb{J}_{\mathcal{F}_n}\|_1 \leq C_2$ is easy since all the coordinates of $\mathbb{J}_{\mathcal{F}_n}$ are uniformly bounded w.r.t. $\boldsymbol{\tau}$ thanks to the upper bound (12) on s' and u' .

Let us now prove that $\|\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})\|_1^{-1} \leq C_3$. Thanks to Corollary 3.6, it suffices to check that $\mathbb{J}_{\mathcal{F}_n}$ is a column-wise (δ, Δ) -M-matrix where δ and Δ depend on the prescribed quantities. First, it is easy to check that

$$(62) \quad j_{LK}^n(\boldsymbol{\tau}) \leq 0 \text{ if } L \neq K, \quad \sum_{L \in \mathcal{T}} j_{LK}^n(\boldsymbol{\tau}) \geq 0, \quad \text{and} \quad 0 \leq \delta_1 \leq j_{KK}^n(\boldsymbol{\tau}) \leq \Delta,$$

where

$$\begin{aligned} \delta &= \alpha_\star \min \left\{ 1; \Delta t^n \min_K \left(\frac{\min_{\sigma \in \mathcal{E}_K} A_\sigma}{m_K} \right) \right\}, \\ \Delta &= \alpha^\star \max_{K \in \mathcal{T}} \left(1 + \frac{\Delta t^n}{m_K} \sum_{\sigma \in \mathcal{E}_K} (m_\sigma g_{K,\sigma}^+ \|\lambda'\|_\infty + A_\sigma) \right). \end{aligned}$$

Therefore, Condition (i) in Definition 3.4 is fulfilled.

It follows from the non-degeneracy condition (11) that for all $\boldsymbol{\tau} \in \mathbb{X}_{\mathcal{T}, \text{int}}$ and all $K \in \mathcal{T}$, either $s'(\tau_K) \geq \alpha_\star$ or $u'(\tau_K) \geq \alpha_\star$. Let K be such that $s'(\tau_K) \geq \alpha_\star$, then

$$\sum_{L \in \mathcal{T}} j_{LK}^n(\boldsymbol{\tau}) \geq \alpha_\star \geq \delta,$$

whence $K \in \mathcal{I}_\delta(\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})^T)$. On the other hand, if $u'(\tau_K) \geq \alpha_\star$ and $\mathcal{E}_{K, \text{ext}} \neq \emptyset$, then

$$\sum_{L \in \mathcal{T}} j_{LK}^n(\boldsymbol{\tau}) \geq \alpha_\star \frac{\Delta t^n}{m_K} \sum_{\sigma \in \mathcal{E}_{K, \text{ext}}} A_\sigma \geq \delta.$$

As a consequence, if $K \notin \mathcal{I}_\delta(\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})^T)$, one has necessarily that $u'(\tau_K) \geq \alpha_\star$ and $\mathcal{E}_{K, \text{ext}} = \emptyset$. But in this case,

$$j_{LK}^n(\boldsymbol{\tau}) \leq -\alpha_\star \frac{\Delta t^n A_\sigma}{m_K} \leq -\delta \quad \text{if } \sigma = K|L \in \mathcal{E}_{K, \text{int}}.$$

The matrix $\mathbb{D}^n \in \mathbb{R}^{\#\mathcal{T} \times \#\mathcal{T}}$ defined by

$$\mathbb{D}_{KK}^n = \alpha_\star \frac{\Delta t^n}{m_K} \sum_{\sigma \in \mathcal{E}_K} A_\sigma, \quad \mathbb{D}_{KL}^n = \begin{cases} -\alpha_\star \frac{\Delta t^n}{m_K} A_\sigma & \text{if } \sigma = K|L \\ 0 & \text{otherwise} \end{cases}$$

is irreducible and admits δ -transmissive paths from K to L for all $(K, L) \in \mathcal{T}^2$. This ensures that $\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})$ admits a δ -transmissive path from any cell $K \in \mathcal{I}_\delta(\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})^T)^c$

to any cell $L \in \mathcal{I}_\delta(\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})^T)$. Therefore, Condition (ii) of Definition 3.4 is fulfilled. $\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})^T$ is then a row-wise (δ, Δ) -M-matrix, thus $\mathbb{J}_{\mathcal{F}_n}(\boldsymbol{\tau})$ is a column-wise (δ, Δ) -M-matrix. \square

The following corollary is a straightforward compilation of Newton-Kantorovich theorem 3.1 with Proposition 3.7.

Corollary 3.8 (local convergence of Newton's method). *There exists $\rho^n > 0$ such that the Newton method (53) converges as soon as $\|\boldsymbol{\tau}^{n,0} - \boldsymbol{\tau}^n\| \leq \rho^n$.*

Remark 3.9 (global convergence for small enough time steps). *The radius ρ^n appearing in Corollary 3.8 can be estimated thanks to [29] as soon as the second order derivatives s'' and u'' of s and u are uniformly bounded. It appears in a fairly natural way that ρ^n is a non-decreasing function of Δt^n . Then choosing Δt^n small enough, the variation between the time steps t^{n-1} and t^n is small and*

$$\|\boldsymbol{\tau}^{n-1} - \boldsymbol{\tau}^n\|_1 = \|\boldsymbol{\tau}^{n,0} - \boldsymbol{\tau}^n\|_1 \leq \rho^n.$$

Therefore, the convergence of the Newton method is ensured if one uses an adaptive time-step algorithm (as for instance in [10]).

3.3. Control of the error induced by the inexact Newton procedure. Assume that $s(\boldsymbol{\tau}^{n-1})$ is exactly known, then the exact solution $\boldsymbol{\tau}^n$ is obtained as the limit of $(\boldsymbol{\tau}^{n,k})_k$. Computing the exact value of $\boldsymbol{\tau}^n$ is impossible and a convenient criterion must be adopted in order to stop the iterative procedure. This yields errors that accumulate along time. The goal of this section is to quantify the error induced by the inexact resolution of the nonlinear system.

3.3.1. One step error estimates. In this section, we assume that $s(\boldsymbol{\tau}^{n-1})$ is exact, and we want to quantify the error corresponding to one single iteration. In what follows, we consider the following residual based stopping criterion:

$$(63) \quad \text{stop the iterative procedure (53) if } \|\mathcal{F}_n(\boldsymbol{\tau}^{n,k})\|_1 = \sum_{K \in \mathcal{T}} |f_K(\tau_K^{n,k})| \leq \epsilon \Delta t^n$$

for some prescribed tolerance $\epsilon > 0$. We denote by $\boldsymbol{\tau}_\epsilon^n = \boldsymbol{\tau}^{n,k}$ when (63) is fulfilled and the loop is stopped. Then the non-degeneracy (58) of $[\mathbb{J}_{\mathcal{F}_n}]^{-1}$ provides directly the following error estimate:

$$\|\boldsymbol{\tau}_\epsilon^n - \boldsymbol{\tau}^n\|_1 \leq C_3 \epsilon \Delta t^n.$$

More than in the variable τ , whose physical sense is unclear, we are interested in evaluating the error on the reconstructed saturation profile. Thanks to the Lipschitz continuity of s —recall the non-degeneracy assumption (11)—, we have

$$(64) \quad \|\pi_{\mathcal{T}} s(\boldsymbol{\tau}_\epsilon^n) - \pi_{\mathcal{T}} s(\boldsymbol{\tau}^n)\|_{L^1(\Omega)} \leq C_4 \epsilon \Delta t^n,$$

where $C_4 = \alpha^* (\max_K m_K) C_3$.

3.3.2. Quantification of the error accumulation. In the previous paragraph, it was assumed that $s(\boldsymbol{\tau}^{n-1})$ was exactly known. In practice, one may consider that we know $s(\boldsymbol{\tau}^0)$ exactly, but merely the approximation of $s(\boldsymbol{\tau}_\epsilon^n)$ obtained after stopping the Newton iterative procedure after a finite number of iterations. Denote by $(s(\boldsymbol{\tau}^n))_{1 \leq n \leq N}$ the iterated exact solutions to the scheme (31)–(32), and by

$(s(\tau_\epsilon^n))_{1 \leq n \leq N}$ the iterated inexact solutions obtained *via* the Newton method with the stopping criterion (63). Define $\mathcal{F}_{n,\epsilon} : \mathbb{X}_{\mathcal{T},\text{int}} \rightarrow \mathbb{X}_{\mathcal{T},\text{int}}$ by

$$(65) \quad \mathcal{F}_{n,\epsilon}(\tau) = \mathcal{F}_n(\tau) + s(\tau^{n-1}) - s(\tau_\epsilon^{n-1}).$$

In particular, one has $\mathbb{J}_{\mathcal{F}_{n,\epsilon}}(\tau) = \mathbb{J}_{\mathcal{F}_n}(\tau)$ for all τ and the results of §3.2 still hold for $\mathcal{F}_{n,\epsilon}$ instead of \mathcal{F}_n . The inexact Newton method then writes

- (1) **Initialization:** define $\tau_\epsilon^0 = \tau^0$ by (28).
- (2) **From t^{n-1} to t^n :**
 - (a) set $\tau_\epsilon^n = \tau_\epsilon^{n-1}$
 - (b) iterate Newton's algorithm until $\|\mathcal{F}_{n,\epsilon}(\tau_\epsilon^n)\|_1 \leq \epsilon \Delta t^n$

Then, as claimed by the following statement, the $L^1(\Omega)$ error on the reconstructed saturation growth at most linearly with time. In particular, no exponential amplification of the error occurs in this context.

Proposition 3.10. *Let $(\tau^n)_n$ be the exact solution to the scheme (31)–(32), and let $(\tau_\epsilon^n)_n$ be the approximate solution computed by the inexact Newton method, then*

$$(66) \quad \|\pi_{\mathcal{T}}s(\tau_\epsilon^n) - \pi_{\mathcal{T}}s(\tau^n)\|_{L^1(\Omega)} \leq C_4 \epsilon t^n, \quad \forall n \in \{0, \dots, N\},$$

where C_4 is introduced as (64).

Proof. We perform the proof by induction. Estimate (66) clearly holds for $n = 0$ since the initialization is exact. Now assume that it holds for $n - 1$. Denote by $\tilde{\tau}_\epsilon^n$ the exact solution of the system $\mathcal{F}_{n,\epsilon}(\tilde{\tau}_\epsilon^n) = \mathbf{0}$. Then it follows from the discussion carried out in §3.3.1 that

$$\|\pi_{\mathcal{T}}s(\tau_\epsilon^n) - \pi_{\mathcal{T}}s(\tilde{\tau}_\epsilon^n)\|_{L^1(\Omega)} \leq C_4 \epsilon \Delta t^n.$$

On the other hand, applying Lemma 2.5, one gets that

$$\|\pi_{\mathcal{T}}s(\tilde{\tau}_\epsilon^n) - \pi_{\mathcal{T}}s(\tau^n)\|_{L^1(\Omega)} \leq \|\pi_{\mathcal{T}}s(\tilde{\tau}_\epsilon^{n-1}) - \pi_{\mathcal{T}}s(\tau^{n-1})\|_{L^1(\Omega)} \leq C_4 \epsilon t^{n-1}.$$

One concludes thanks to the triangle inequality. \square

4. NUMERICAL VALIDATION OF THE APPROACH

In the previous section, we have shown that as long as the parametrization $\tau \mapsto (s(\tau), u(\tau))$ satisfies the condition (11), the Newton's method applied to the problem (52) exhibits local convergence. Moreover, due to Proposition 3.10, the error resulting from an inexact Newton's method can be efficiently controlled. Remark that the constant C_4 in (66) depends on the nonlinearities S and η only through the quantities α_\star and α^\star . Therefore, the estimate (66) is robust with respect to the hydrodynamic properties of the soil. We illustrate this fact by numerical experiments presented below. As an example of mobility/capillary pressure relations we consider a popular Brooks-Corey model [6] for which we compare the efficiency of Newton's method resulting from the parametrization $u(\tau) = \tau$ and the one satisfying (11) with $\alpha_\star = \alpha^\star = 1$.

Let $p_b < 0$ and $\beta > 0$, in Brooks-Corey model the saturation and the mobility functions are given by

$$S(p) = \begin{cases} \left(\frac{p}{p_b}\right)^{-\beta} & \text{if } p < p_b, \\ 1 & \text{if } p \geq p_b, \end{cases}$$

and

$$\lambda(s) = s^{3+\frac{2}{\beta}},$$

providing in terms of Kirchhoff transform

$$\tilde{S}(u) = \begin{cases} \left(\frac{u}{u_b}\right)^{\frac{1}{\eta}} & \text{if } u < u_b, \\ 1 & \text{if } u \geq u_b, \end{cases}$$

with $\eta = \beta + 3 + \frac{1}{\beta}$ and $u_b = -\frac{p_b}{\beta\eta}$.

Remark that the parametrization based on $u(\tau) = \tau$ and $s(\tau) = \tilde{S}(u(\tau))$, referred as u -formulation, do not satisfy (11) since the derivative of $\tilde{S}(u)$ is singular at $u = 0$. As an alternative we consider the parametrization $\tau \mapsto (s(\tau), u(\tau))$ defined by the equation $\max(s'(\tau), u'(\tau)) = 1$ and the condition $s(0) = 0$, to which we refer as τ -formulation. We obtain the following explicit formulas

$$s(\tau) = \begin{cases} \tau & \text{if } \tau < \tau_*, \\ S(\tau - \tau_* + u_b \tau_*^\eta) & \text{if } \tau \geq \tau_*, \end{cases}$$

and

$$u(\tau) = \begin{cases} u_b \tau^\eta & \text{if } \tau < \tau_*, \\ \tau - \tau_* + u_b \tau_*^\eta & \text{if } \tau \geq \tau_*, \end{cases}$$

with $\tau_* = \min\left((\eta u_b)^{\frac{1}{1-\eta}}, 1\right)$.

4.1. First test case. We consider a bidimensional porous domain $\Omega = (0, 1) \times (0, 1)$, which is initially very dry with $s_0(\mathbf{x}) = 10^{-6}$ in Ω . The water is injected at the pressure $p_D = 1$ through the portion of an upper boundary $\Gamma_D = \{(x_1, x_2) \mid x_1 \in (0, 0.3), x_2 = 1\}$. The gravity vector is given by $\mathbf{g} = -\nabla x_2$ and a zero flux boundary condition is prescribed on $\partial\Omega \setminus \Gamma_D$. The computations are performed on two quasi-uniform space discretizations of Ω composed of 396 and 1521 Voronoï cells referred as Mesh 1 and Mesh 2. The final time is set to $T = 0.7$ and the time step is equal to 0.01. Figure 2 shows, for $\beta = 4$ and $p_b = -10^{-2}$, the distribution of saturation and generalized pressure u at different times. The results are visualized on triangular Delaunay mesh dual to Mesh 1.

In order to challenge the robustness of both formulations, we set $p_b = -10^{-2}$ and we let the parameter β take values in the set $\{1, 2, 4, 8, 16\}$. For each value of β we compute, using τ -formulation and tolerance $\epsilon_{ref} = 10^{-12}$, the reference solution denoted by $(\boldsymbol{\tau}_\beta^n)_{n \in \{1, \dots, N\}} \in \mathbb{X}_D$. Then, for both formulations and for the values of $\epsilon \in \{10^{-2}, 10^{-4}, 10^{-6}, 10^{-8}, 10^{-10}, 10^{-12}\}$, we perform the calculations measuring the total number of Newton's iteration and the deviation, in the discrete $L^\infty(L^1)$ norm, of the “observable” variables u and s from the reference solution.

For a given value of β and of the tolerance ϵ , we denote by $(\bar{\boldsymbol{\tau}}_{\beta, \epsilon}^n)_{n \in \{1, \dots, N\}} \in \mathbb{X}_D$ and $(\boldsymbol{\tau}_{\beta, \epsilon}^n)_{n \in \{1, \dots, N\}} \in X_D$ the approximate solution of (52) obtained using the u -formulation and τ -formulation respectively. The error produced by inexact Newton's method is measured by the quantities

$$\overline{err}_{\beta, \epsilon}^u = \frac{\|\pi_D \bar{\boldsymbol{\tau}}_{\beta, \epsilon}^n - \pi_D u(\boldsymbol{\tau}_\beta^n)\|_{L^\infty(0, T; L^1(\Omega))}}{\|\pi_D u(\boldsymbol{\tau}_\beta^n)\|_{L^\infty(0, T; L^1(\Omega))}}$$

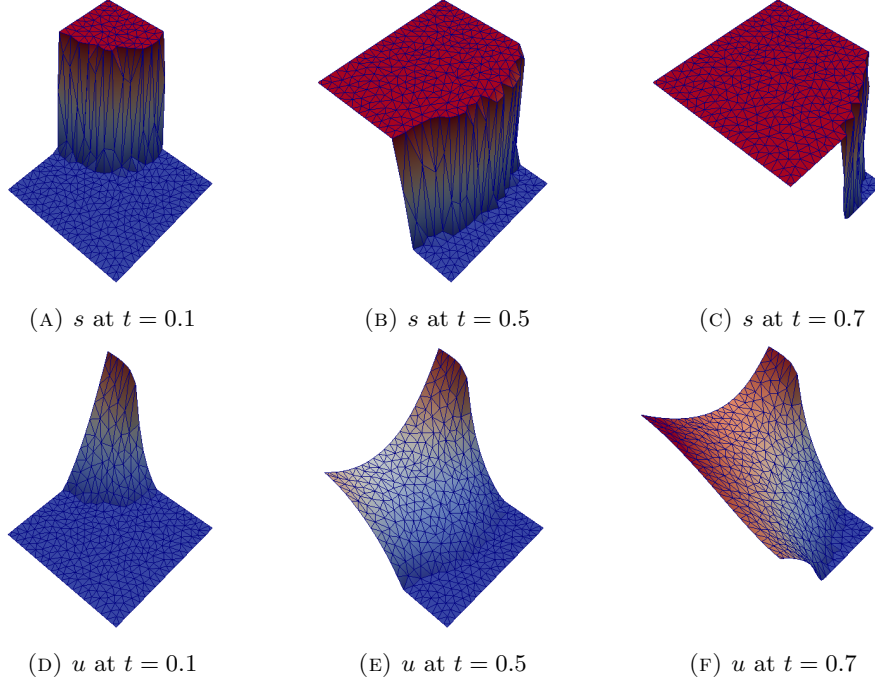


FIGURE 2. Snapshots of the reference solution at different times.

and

$$\overline{err}_{\beta,\epsilon}^s = \frac{\|\pi_{\mathcal{D}}\tilde{S}(\bar{\tau}_{\beta,\epsilon}^n) - \pi_{\mathcal{D}}s(\tau_{\beta}^n)\|_{L^\infty(0,T;L^1(\Omega))}}{\|\pi_{\mathcal{D}}s(\tau_{\beta}^n)\|_{L^\infty(0,T;L^1(\Omega))}}$$

for u -formulation, and

$$err_{\beta,\epsilon}^u = \frac{\|\pi_{\mathcal{D}}u(\tau_{\beta,\epsilon}^n) - \pi_{\mathcal{D}}u(\tau_{\beta}^n)\|_{L^\infty(0,T;L^1(\Omega))}}{\|\pi_{\mathcal{D}}u(\tau_{\beta}^n)\|_{L^\infty(0,T;L^1(\Omega))}}$$

and

$$err_{\beta,\epsilon}^s = \frac{\|\pi_{\mathcal{D}}s(\tau_{\beta,\epsilon}^n) - \pi_{\mathcal{D}}s(\tau_{\beta}^n)\|_{L^\infty(0,T;L^1(\Omega))}}{\|\pi_{\mathcal{D}}s(\tau_{\beta}^n)\|_{L^\infty(0,T;L^1(\Omega))}}$$

for τ -formulation.

Figure 3 exhibits, for Meshes 1 and 2, the behavior of the relative error $err_{\beta,\epsilon}^\xi$ and $\overline{err}_{\beta,\epsilon}^\xi$, $\xi = u, s$ as the function of an average number of iterations per time step required by Newton's method in order to converge up to the given tolerance ϵ . We observe that in order to achieve the same precision u -formulation require a much larger number of iterations then τ -formulation. Moreover the number of Newton's iterations, for u -formulation, increases with β , whereas τ -formulation remains robust with respect to this parameter. The contrast in the efficiency of two formulations is amplified as the mesh is refined.

4.2. Second test case. The goal of this test case is to give a numerical evidence that the inexact Newton's method applied to the u -formulation produces large

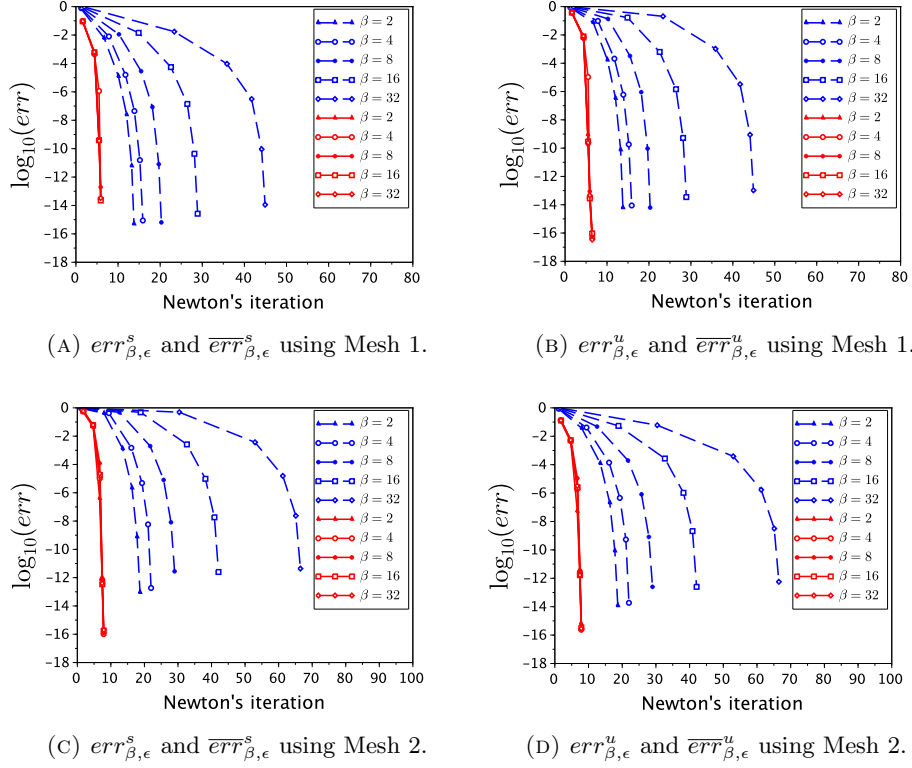


FIGURE 3. Relative error as the function of the average number of Newton's iterations per time step using τ -formulation (red solid lines) and u -formulation (blue dashed lines).

errors due to troubles in the conservation of mass. To do so, we prescribe a zero-flux condition on the whole boundary of $\Omega = (0, 1) \times (0, 1)$, whereas the initial saturation satisfies

$$s_0 = \begin{cases} 0.5 & \text{in } \Omega', \\ 10^{-6} & \text{in } \Omega \setminus \Omega' \end{cases}$$

with $\Omega' = \{(x_1, x_2) \mid x_1 < 0.5 \text{ and } x_2 > 0.5\}$. We set $\beta = 4$ and $p_b = -10^{-2}$. The effects of gravity neglected so that the flow is only driven by diffusion. Figure 5 exhibits, for different times, the saturation field associated with the reference solution, which is computed using τ -formulation and the tolerance $\epsilon_{ref} = 10^{-16}$. Since the flow is very slow the large time steps are needed, we set $T = 10^5$ and $\Delta t = 10^3$. Computations are performed using Mesh 1. Let $M = \int_{\Omega} s_0 d\mathbf{x}$, we define the relative mass conservation error by

$$\overline{mass\ err}_{\beta,\epsilon}^s = \frac{1}{M} \max_{n \in \{1, \dots, N\}} \left| \int_{\Omega} \pi_{\mathcal{D}} \tilde{S}(\bar{\tau}_{\beta,\epsilon}^n) d\mathbf{x} - M \right|$$

and

$$mass\ err_{\beta,\epsilon}^s = \frac{1}{M} \max_{n \in \{1, \dots, N\}} \left| \int_{\Omega} \pi_{\mathcal{D}} s(\tau_{\beta,\epsilon}^n) d\mathbf{x} - M \right|.$$

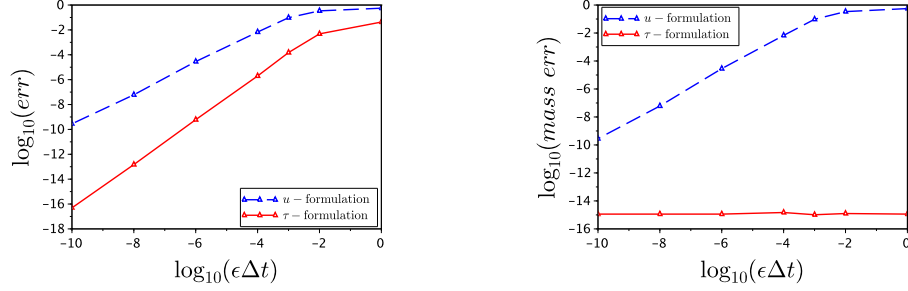


FIGURE 4. At left $err_{\beta,\epsilon}^s$ (red solid lines) and $\overline{err}_{\beta,\epsilon}^s$ (blue dashed lines) as the function of ϵ , at right the relative mass conservation error as the function of ϵ .

Figure 4 exhibits the $L^\infty(L^1)$ relative saturation error $err_{\beta,\epsilon}^s$, $\overline{err}_{\beta,\epsilon}^s$, and the relative mass conservation error $mass\ err_{\beta,\epsilon}^s$ and $\overline{mass\ err}_{\beta,\epsilon}^s$ as the functions of ϵ . As one can see the $L^\infty(L^1)$ error produced by u -formulation is dominated by the mass conservation error. Remark that even for the rather small values $\epsilon = 10^{-6}$ or $\epsilon = 10^{-7}$, the error produced by u -formulation is still significant (see Figures 6). In contrast τ -formulation leads to much smaller errors and, for any value of ϵ , conserves the mass up to a precision of order 10^{-15} .

The very high accuracy for mass conservation observed with the τ -formulation can be explained as follows. With the values of the parameters β, p_b we have chosen, one has $s(\tau) = \tau$ for $\tau \in [0, 1]$ and hence for all $\tau \in \Omega \times (0, T)$ in view of initial and boundary conditions. In addition, in view of (32) we have

$$\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}(\tau^n) = 0, \quad \forall \tau^n \in \mathbb{X}_D.$$

Therefore, at each step of inexact Newton's method, the flux contribution globally offset, hence we have

$$\sum_{K \in \mathcal{T}} m_K \left(s(\tau_{\epsilon,K}^{n,k}) + s'(\tau_{\epsilon,K}^{n,k})(\tau_{\epsilon,K}^{n,k+1} - \tau_{\epsilon,K}^{n,k}) - s(\tau_{\epsilon,K}^{n-1}) \right) = 0.$$

Since $\tau \mapsto s(\tau)$ is linear, one has $s(\tau_{\epsilon,K}^{n,k+1}) = s(\tau_{\epsilon,K}^{n,k}) + s'(\tau_{\epsilon,K}^{n,k})(\tau_{\epsilon,K}^{n,k+1} - \tau_{\epsilon,K}^{n,k})$, which implies that the mass is exactly conserved (assuming that linear algebraic computations are exact) at each iteration of Newton's method.

REFERENCES

- [1] H. W. Alt and S. Luckhaus. Quasilinear elliptic-parabolic differential equations. *Math. Z.*, 183(3):311–341, 1983.
- [2] B. Andreianov, C. Cancès, and A. Moussa. A nonlinear time compactness result and applications to discretization of degenerate parabolic-elliptic PDEs. HAL: hal-01142499, 2015.
- [3] J. Bear and Y. Bachmat. *Introduction to modeling of transport phenomena in porous media*. Kluwer Academic Publishers, Dordrecht, The Netherlands, 1990.
- [4] L. Bergamaschi and M. Putti. Mixed finite elements and Newton-type linearizations for the solution of Richards' equation. *Int. J. Numer. Meth. Eng.*, 45(8):1025–1046, 1999.

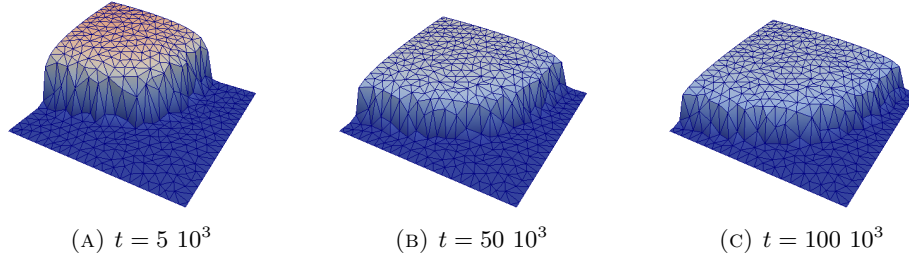
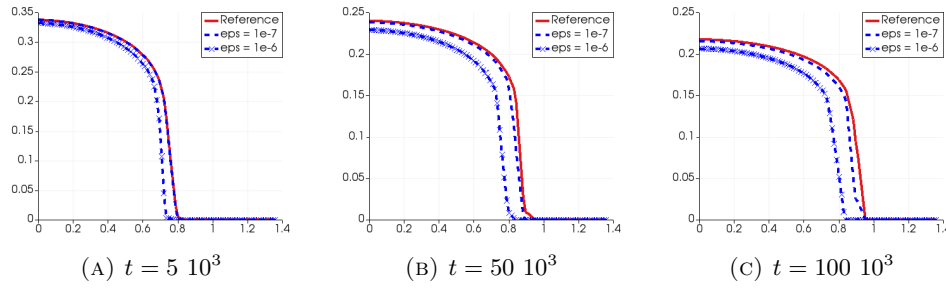


FIGURE 5. Saturation field of the reference solution at different times.

FIGURE 6. Saturation profile of the reference solution along the line $y = 1 - x$ compared to the saturation profiles of the approximate solutions computed using u -formulation and $\epsilon = 10^{-8}$ and 10^{-9} .

- [5] K. Brenner, M. Groza, L. Jeannin, R. Masson, and J. Pellerin. Immiscible two-phase darcy flow model accounting for vanishing and discontinuous capillary pressures: application to the flow in fractured porous media. *Proceedings of ECMOR XV*, 2016. To appear.
- [6] R. H. Brooks and A. T. Corey. Hydraulic properties of porous media and their relation to drainage design. *Transactions of the ASAE*, 7(1):0026–0028, 1964.
- [7] X.-C. Cai and D. E. Keyes. Nonlinearly preconditioned inexact Newton algorithms. *SIAM J. Sci. Comp.*, 24(1):183–200, 2002.
- [8] C. Cancès and T. Gallouët. On the time continuity of entropy solutions. *J. Evol. Equ.*, 11(1):43–55, 2011.
- [9] C. Cancès and C. Guichard. Convergence of a nonlinear entropy diminishing Control Volume Finite Element scheme for solving anisotropic degenerate parabolic equations. *Math. Comp.*, 85(298):549–580, 2016.
- [10] C. Cancès and C. Guichard. Numerical analysis of a robust free energy-diminishing Finite Volume scheme for degenerate parabolic equations with gradient structure. preprint HAL: hal-01119735, accepted for publication in Found. Comput. Math., 2016.
- [11] C. Cancès and M. Pierre. An existence result for multidimensional immiscible two-phase flows with discontinuous capillary pressure field. *SIAM J. Math. Anal.*, 44(2):966–992, 2012.
- [12] J. Carrillo. On the uniqueness of the solution of the evolution dam problem. *Nonlinear Anal.*, 22(5):573–607, 1994.
- [13] J. Carrillo. Entropy solutions for nonlinear degenerate problems. *Arch. Ration. Mech. Anal.*, 147(4):269–361, 1999.
- [14] J. Carrillo. Conservation laws with discontinuous flux functions and boundary condition. *J. Evol. Equ.*, 3(2):283–301, 2003.
- [15] F. H. Clarke. *Optimization and nonsmooth analysis*. Siam, 1990.
- [16] K. Deimling. *Nonlinear functional analysis*. Springer-Verlag, Berlin, 1985.

- [17] H.-J.G. Diersch and P. Perrochet. On the primary variable switching technique for simulating unsaturated-saturated flows. *Adv. Water Resour.*, 23(3):271–301, 1999.
- [18] V. Dolean, M. J. Gander, W. Kheriji, F. Kwok, and R. Masson. Nonlinear preconditioning: how to use a nonlinear Schwarz method to precondition Newton’s method. preprint HAL:hal-01171167, July 2015.
- [19] J. Droniou and R. Eymard. Uniform-in-time convergence of numerical methods for non-linear degenerate parabolic equations. *Numer. Math.*, 132(4):721–766, 2016.
- [20] J. Droniou, R. Eymard, T. Gallouët, and R. Herbin. Gradient schemes: a generic framework for the discretisation of linear, nonlinear and nonlocal elliptic and parabolic equations. *Math. Models Methods Appl. Sci.*, 23(13):2395–2432, 2013.
- [21] R. Eymard, T. Gallouët, M. Ghilani, and R. Herbin. Error estimates for the approximate solutions of a nonlinear hyperbolic equation given by finite volume schemes. *IMA J. Numer. Anal.*, 18(4):563–594, 1998.
- [22] R. Eymard, T. Gallouët, C. Guichard, R. Herbin, and R. Masson. TP or not TP, that is the question. *Comput. Geosci.*, 18:285–296, 2014.
- [23] R. Eymard, T. Gallouët, and R. Herbin. Finite volume methods. Ciarlet, P. G. (ed.) et al., in Handbook of numerical analysis. North-Holland, Amsterdam, pp. 713–1020, 2000.
- [24] R. Eymard, T. Gallouët, and R. Herbin. Discretization of heterogeneous and anisotropic diffusion problems on general nonconforming meshes sushi: a scheme using stabilization and hybrid interfaces. *IMA J. Numer. Anal.*, 30(4):1009–1043, 2010.
- [25] R. Eymard, M. Gutnic, and D. Hilhorst. The finite volume method for Richards equation. *Comput. Geosci.*, 3(3-4):259–294, 1999.
- [26] J. Fuhrmann. Existence and uniqueness of solutions of certain systems of algebraic equations with off-diagonal nonlinearity. *Appl. Numer. Math.*, 37:359–370, 2001.
- [27] J. Fuhrmann and H. Langmach. Stability and existence of solutions of time-implicit finite volume schemes for viscous nonlinear conservation laws. *Appl. Numer. Math.*, 37:201–230, 2001.
- [28] G. Gagneux and M. Madaune-Tort. Unicité des solutions faibles d’équations de diffusion-convection. *C. R. Acad. Sci. Paris Sér. I Math.*, 318(10):919–924, 1994.
- [29] W. B. Gragg and R. A. Tapia. Optimal error bounds for the Newton-Kantorovich theorem. *SIAM J. Numer. Anal.*, 11:10–13, 1974.
- [30] W. Jäger and J. Kačur. Solution of porous medium type systems by linear approximation schemes. *Numer. Math.*, 60(3):407–427, 1991.
- [31] W. Jäger and J. Kačur. Solution of doubly nonlinear and degenerate parabolic problems by relaxation schemes. *RAIRO Modél. Math. Anal. Numér.*, 29(5):605–627, 1995.
- [32] L. V. Kantorovich. On Newtons method for functional equations. *Dokl. Akad. Nauk SSSR*, 59(7):1237–1240, 1948.
- [33] F. Lehmann and P. H. Ackerer. Comparison of iterative methods for improved solutions of the fluid flow equation in partially saturated porous media. *Transport in Porous Media*, 31(3):275–292, 1998.
- [34] J. Leray and J. Schauder. Topologie et équations fonctionnelles. *Ann. Sci. École Norm. Sup.* (3), 51:45–78, 1934.
- [35] F. List and F. A. Radu. A study on iterative methods for solving Richards equation. *Comput. Geosci.*, online first:1–13, 2016.
- [36] S. Martin and J. Vovelle. Convergence of implicit finite volume methods for scalar conservation laws with discontinuous flux function. *ESAIM Math. Model. Numer. Anal.*, 42(5):699728, 2008.
- [37] J. M. Ortega. The Newton-Kantorovich theorem. *Amer. Math. Monthly*, 75:658–660, 1968.
- [38] J. M. Ortega and W. C. Rheinboldt. *Iterative solution of nonlinear equations in several variables*. Academic Press, New York-London, 1970.
- [39] F. Otto. L^1 -contraction and uniqueness for quasilinear elliptic-parabolic equations. *J. Differential Equations*, 131:20–38, 1996.
- [40] I. S. Pop, F. Radu, and P. Knabner. Mixed finite elements for the Richards equation: linearization procedure. *J. Comput. Appl. Math.*, 168(1):365–373, 2004.
- [41] L. Qi and J. Sun. A nonsmooth version of Newton’s method. *Mathematical programming*, 58:353–367, 1993.

- [42] F. A. Radu, I. S. Pop, and P. Knabner. Newton-type methods for the mixed finite element discretization of some degenerate parabolic equations. In *Numerical mathematics and advanced applications*, pages 1192–1200. Springer, 2006.
- [43] L. A. Richards. Capillary conduction of liquids through porous mediums. *Journal of Applied Physics*, 1(5):318–333, 1931.
- [44] X. Wang and H. A. Tchelepi. Trust-region based solver for nonlinear transport in heterogeneous porous media. *J. Comput. Phys.*, 253:114–137, 2013.
- [45] R. Younis, H. A. Tchelepi, and K. Aziz. Adaptively localized continuation-Newton method–nonlinear solvers that converge all the time. *SPE Journal*, 15(02):526–544, 2010.
- [46] R. L. Zarba, E. T. Bouloutas, and M. Celia. General mass-conservative numerical solution for the unsaturated flow equation. *Water Resour. Res.*, 26(7):1483–1496, 1990.

Konstantin BRENNER

Laboratoire Jean-Alexandre Dieudonné, Université de Nice Sophia Antipolis,
Team Coffee INRIA Sophia Antipolis Méditerranée,
06108 Nice Cedex 02, France.

konstantin.brenner@unice.fr

Clément CANCES

Team Rapsodi INRIA Lille - Nord Europe,
40, avenue Halley, 59650 Villeneuve d’Ascq, France.

clement.cances@inria.fr